# **Type-Directed Continuation Allocation\***

Zhong Shao and Valery Trifonov

Dept. of Computer Science Yale University New Haven, CT 06520-8285 {shao,trifonov}@cs.yale.edu

**Abstract.** Suppose we translate two different source languages,  $L_1$  and  $L_2$ , into the same intermediate language; can they safely interoperate in the same address space and under the same runtime system? If  $L_1$  supports first-class continuations (call/cc) and  $L_2$  does not, can  $L_2$  programs call arbitrary  $L_1$  functions? Would the fact of possibly calling  $L_1$  impose restrictions on the implementation strategy of  $L_2$ ? Can we compile  $L_1$  functions that do not invoke call/cc using more efficient techniques borrowed from the  $L_2$  implementation? Our view is that the implementation of a common intermediate language ought to support the so-called pay-as-you-go efficiency: first-order monomorphic functions should be compiled as efficiently as in C and assembly languages, even though they may be passed to arbitrary polymorphic functions that support advanced control primitives (e.g. call/cc). In this paper, we present a typed intermediate language with effect and resource annotations, ensuring the safety of inter-language calls while allowing the compiler to choose continuation allocation strategies.

# 1 Introduction

Safe interoperability requires resolving a host of issues including mixed data representations, multiple function calling conventions, and different implementation protocols. Existing approaches to language interoperability either separate code written in different languages into different address spaces or have the unsafe, ad hoc and insecure foreign function call interface.

We position our further discussion of language interoperability in the context of a system hosting multiple languages, each safe in isolation. The supported languages may range from first-order monomorphic (e.g. a safe subset of C, or safe-C for short) to higher-order languages with advanced control, e.g. ML with first-class continuations. We assume that all languages have type systems which ensure runtime safety of accepted programs. In other words, in this paper we do not attempt to solve the problem of cooperating safely with programs written in unsafe languages, which in general can

<sup>\*</sup> This research was sponsored in part by the DARPA ITO under the title "Software Evolution using HOT Language Technology", DARPA Order No. D888, issued under Contract No. F30602-96-2-0232, and in part by an NSF CAREER Award CCR-9501624, and NSF Grant CCR-9633390. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government.

maintaining the data needed, and of including suggestions for reducing	lection of information is estimated to completing and reviewing the collect this burden, to Washington Headqu uld be aware that notwithstanding ar OMB control number.	ion of information. Send comments arters Services, Directorate for Information	regarding this burden estimate mation Operations and Reports	or any other aspect of the , 1215 Jefferson Davis	is collection of information, Highway, Suite 1204, Arlington
1. REPORT DATE <b>2005</b>		2. REPORT TYPE		3. DATES COVE	RED
4. TITLE AND SUBTITLE		5a. CONTRACT NUMBER			
Type-Directed Cor		5b. GRANT NUMBER			
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)  Defense Advanced Research Projects Agency,3701 North Fairfax  Dr,Arlington,VA,22203-1714				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAIL  Approved for publ	LABILITY STATEMENT ic release; distributi	ion unlimited			
13. SUPPLEMENTARY NO	OTES				
14. ABSTRACT see report					
15. SUBJECT TERMS					
16. SECURITY CLASSIFIC	17. LIMITATION OF	18. NUMBER	19a. NAME OF		
a. REPORT <b>unclassified</b>	b. ABSTRACT <b>unclassified</b>	c. THIS PAGE unclassified	- ABSTRACT	OF PAGES 20	RESPONSIBLE PERSON

**Report Documentation Page** 

Form Approved OMB No. 0704-0188 only be achieved at the expense of "sandboxing" the unsafe calls or complex and incomplete analyses of the unsafe code.

We believe that interoperability requires a serious and more formal treatment. As a first step, this paper describes a novel type-based technique to support principled language interoperation among languages with different protocols for allocation of activation records. Our framework allows programs written in multiple languages with overlapping features to interact with each other safely and reliably, yet without restricting the expressiveness of each language.

An interoperability scheme for activation record allocation should be

- safe: it should not be possible to violate the runtime safety of a language by calling a foreign function;
- expressive: the scheme should allow inter-language function calls;
- efficient: a language implementation should not be forced to use suboptimal methods for its own features in order to provide support for other languages' features.
   For instance a language that does not use call/cc should not have to be implemented using heap-based allocation of activation records.

Our solution is to ensure safety by using a common typed intermediate language [21] into which all of the source languages are translated. To maintain safety in an expressive interoperability scheme the type system is extended with annotations of the *effects* of the evaluation of a term, e.g. an invocation of call/cc, and polymorphic types with effect variables, allowing a higher-order function to be invoked with arguments coming from languages with different sets of effects. The central novelty of our approach is the introduction of annotations of the *resources* necessary for the realization of the effects of an evaluation; for instance a continuation heap may be required when invoking call/cc. Thus our type system can be used to support implementation efficiency by keeping track of the available language-dependent resources, and safety by allowing semantically correct inter-language function calls but banning semantically incorrect ones. In addition to providing safety, making resource handling explicit also opens new opportunities for code optimization beyond what a foreign function call mechanism can offer.

A common intermediate language like FLINT [20, 21] will likely support a very rich set of features to accommodate multiple source languages. Some of these features may impose implementation restrictions; for example, a practical implementation of first-class continuations (as in SML/NJ or Scheme) often requires the use of advanced stack representations [8] or heap-based activation records [22]. However in some cases stack-based allocation may be more efficient, and ideally we would like to have a compiler that can take advantage of it as long as this does not interfere with the semantic correctness of first-class continuations. Similarly, when compiling a simple safe-C-like language with no advanced control primitives (e.g., call/cc) into FLINT, we may prefer to compile it to code that uses the simple sequential stack of standard C; programs written in ML or Scheme using these safe-C functions must then follow the same allocation strategy when invoking them. This corresponds to the typical case of writing low-level systems modules in C and providing for their use in other languages, therefore we assume this model in the sequel, but the dual problem of compiling safe-C functions

calling arbitrary ML functions by selectively imposing heap allocation on safe-C is similarly represented and solved within our system.

Thus our goal is efficient and expressive interoperability between code fragments written in languages using possibly different allocation disciplines for activation records, for instance, ML with heap allocation and safe-C with stack allocation. The following properties of the interoperability framework are essential for achieving this goal:

- ML and safe-C code should interoperate safely with each other within the same address space.
- All invocations of safe-C functions in ML functions should be allowed (provided they are otherwise type-correct).
- Only the invocations of ML functions that do not capture continuations should be allowed in safe-C functions.
- Any activation record that can potentially be captured as part of a first-class continuation should always be allocated on the heap (or using some fancy stack-chunkbased representations [8]).
- It should be possible to use stack allocation for activation records of ML functions when they are guaranteed not to be captured with a first-class continuation.
- The selection of allocation strategy should be decoupled from the actual function call.

The last property gives the compiler the freedom to switch allocation strategies more efficiently, instead of following a fixed foreign function interface mechanism. For example, an implementation of ML may use heap allocation of activation records by default to provide support for continuation capture. However, in cases when the compiler can prove that a function's activation record is not going to be accessible from any captured continuation, its allocation discipline is ambiguous; stack allocation may be preferred if the function invokes, or is invoked by, safe-C functions which use stack allocation. This specialization of code to a different allocation strategy effectively creates regions of ML code compiled in "safe-C mode" with the aim of avoiding the switch between heap and stack allocation on every cross-language call. In general, the separation of the selection of allocation strategy from the call allows its treatment as a commodity primitive operation and subjects it to other code-motion optimizations, e.g. hoisting it out of loops.

The proposed method can be applied to achieving more efficient interoperability with existing foreign code as well, although obviously in this case the usual friction between safety and efficiency can only be eased but not removed. In particular the possibility to select the allocation strategy switch point remains, thus higher efficiency can still be achieved while satisfying a given safety policy by specializing safe code to "unsafe mode" (e.g. for running with stack allocation within a sand-box).

# 2 A Resourceful Intermediate Language

To satisfy the requirements for efficient interoperability, outlined in the previous section, we define an A-normal-form-based typed intermediate language RL (Figure 1) with types having effect and resource annotations. Intuitively, an effect annotation such

as CC indicates that a computation may capture a continuation by performing call/cc; a resource annotation such as H (continuation heap) or S (continuation stack) means that the corresponding runtime resource must be available to the computation. Non-trivial effects can be primitive, effect variables, or unions of effects; commutativity and associativity of the union with  $\emptyset$  as a unit are consistent with the typing rules and we assume them for brevity of notation. Each effect can only occur when the proper resources are available, e.g. CC would require the use of heap-based activation record allocation. Both the effect and resource usage annotations are inferred during the translation from the source language to the intermediate language, and can be used to assist code generation and to check the validity of cross-language function calls.

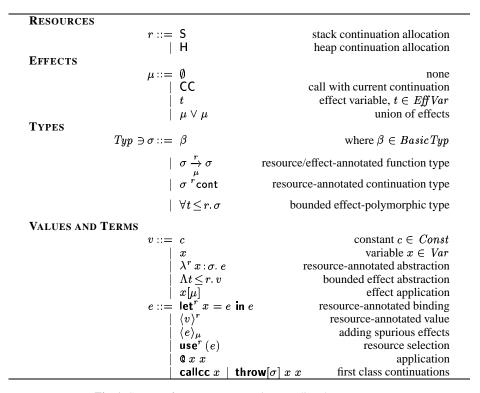


Fig. 1. Syntax of a resource-aware intermediate language RL

The resources required and effects produced by a function are made explicit in its type. A continuation can potentially produce all effects possible with the set of resources available at the point of its capture; for that reason continuation types only have a resource annotation.

<sup>&</sup>lt;sup>1</sup> In this paper, we focus on application of this system to interoperability issues related to continuation allocation, but more diverse sets of resources will be necessary in a realistic language.

Function abstractions are annotated with the resources they may require and will maintain. In a higher-order language the effect of the evaluation of a function application may depend on the effects of its functional arguments; this dependence is expressed by means of effect polymorphism. Polymorphic abstractions introduce variables ranging over the set of possible effects of the term. Since the possible effects are determined by the available resources, we have bounded effect polymorphism; the relation  $\mu \leq r$  (defined in the context of an effect environment in Figure 3) reflects the dependence between effects and resources, e.g. that **callcc** can only be performed if continuations are heap-allocated. The effect application  $x[\mu]$  instantiates the body of the polymorphic abstraction to which x is bound. The language construct  $\mathbf{use}^r$  (e) serves to mark the point where a change in the allocation strategy for activation records is required. Instead of having effect subsumption the language is equipped with a construct  $\langle e \rangle_{\mu}$  for explicitly increasing the set of effects of e to include  $\mu$ .

*Example 1.* The use of resource annotations to select allocation strategies is shown in the *RL* code below which includes extra type annotations for clarity.

$$\begin{split} \textbf{let}^{H} \\ & \text{applyToInt} \\ & = \langle \Lambda t \leq H. \ \lambda^{H} \ f \colon \text{Int} \ \frac{H}{t} \ \text{Int}. \ @ \ f \ 42 \rangle^{H} \\ & : \forall t \leq H. \ (\text{Int} \ \frac{H}{t} \ \text{Int}) \ \frac{H}{t} \ \text{Int} \\ & \text{add1\_CC} \\ & = \langle \lambda^{H} \ x \colon \text{Int}. \\ & \textbf{let}^{H} \ c = \langle \lambda^{H} \ k \colon \text{Int} \ ^{H} \text{cont}. \\ & \textbf{let}^{H} \ z = @ \ \text{succ} \ x \ \textbf{in} \ \textbf{throw}[\text{Int}] \ k \ z \rangle^{H} \\ & \quad \text{in} \ \textbf{callcc} \ c \rangle^{H} \\ & \quad : \text{Int} \ \frac{H}{cC} \ \text{Int} \\ & \text{add1\_Pure} \\ & = \langle \lambda^{S} \ x \colon \text{Int}. \ @ \ \text{succ} \ x \rangle^{H} \\ & \quad : \text{Int} \ \frac{S}{\emptyset} \ \text{Int} \\ & \text{add1\_Wrapped} \\ & = \langle \lambda^{H} \ x \colon \text{Int}. \ \textbf{use}^{S} \ (@ \ \text{add1\_Pure} \ x) \rangle^{H} \\ & \quad : \text{Int} \ \frac{H}{\emptyset} \ \text{Int} \\ & \quad \text{in} \ @ \ (\text{applyToInt}[CC]) \ \text{add1\_CC} \ ; \\ & \quad @ \ (\text{applyToInt}[\emptyset]) \ \text{add1\_Wrapped} \\ \end{split}$$

The function applyToInt is polymorphic in the effect of its parameter, but the parameter's resource requirements are fixed – it must use heap allocation. We consider two applications of applyToInt. The argument in the first, add1\_CC, is a function invoking callcc, which consequently uses heap allocation; on the other hand the argument in the second application, add1\_Pure, is pure and uses stack allocation. It is therefore incorrect to apply applyToInt to add1\_Pure. We use a wrapper to coerce it to the proper type:

we apply applyToInt to add1\_Wrapped whose activation record is heap-allocated, and whose function is to switch to stack allocation (via **use**<sup>S</sup>) before calling add1\_Pure. Heap allocation is resumed upon return from add1\_Pure.

# 3 Two Source Languages

To further illustrate the advantages of this system we consider the problem of translating into RL two source languages (Figure 2): a language HL with control operators (**callcc** and **throw**), implemented using heap-based allocation of activation records, and a language SL which always uses stack allocation. HL also allows declaring at the top of a program the identifiers of entities imported from SL code. The type systems of these languages are assumed monomorphic for simplicity, since polymorphism in types is largely orthogonal to the effect polymorphism of RL.

```
SL \text{ TYPES} \qquad \qquad \tau_{SL} ::= \beta \mid \tau_{SL} \rightarrow \tau_{SL} \\ SL \text{ TERMS} \qquad \qquad e_{SL} ::= c \mid x \mid \lambda x : \tau_{SL} \cdot e_{SL} \mid e_{SL} e_{SL} \mid \text{ let } x = e_{SL} \text{ in } e_{SL} \\ HL \text{ TYPES} \qquad \qquad \tau_{HL} ::= \beta \mid \tau_{HL} \rightarrow \tau_{HL} \mid \tau_{HL} \text{ cont} \\ HL \text{ TERMS} \qquad \qquad e_{HL} ::= c \mid x \mid \lambda x : \tau_{HL} \cdot e_{HL} \mid e_{HL} e_{HL} \mid \text{ let } x = e_{HL} \text{ in } e_{HL} \\ \qquad \qquad \qquad \mid \text{ callcc } e_{HL} \mid \text{ throw}[\tau_{HL}] e_{HL} e_{HL} \\ HL \text{ PROGRAMS} \qquad p_{HL} ::= e_{HL} \mid \text{ external}(SL) \ x : \tau_{SL} \text{ in } p_{HL} \end{aligned}
```

Fig. 2. Syntax of the source languages SL and HL

The resource annotations in RL provide information about handling of the stack and heap resources, necessary in the following situations:

- when calling from HL a function written in SL, which may require switching from heap allocation of activation records to allocation on the stack used by SL; the heap resource must be preserved for use upon return from SL code.
- when calling an HL function from SL code, which is only semantically sound when the evaluation of the function does not capture a continuation, since part of the continuation data is stack-allocated; the type system maintains information about the possible effects of the evaluation, in this case whether callcc might be invoked.
- when selecting an allocation strategy for HL functions called (directly or indirectly) from within SL code; either their activation records must be allocated on the SL stack, or the latter must be preserved and restored upon return to SL.
- when selecting an allocation strategy for HL code invoking SL functions but not callcc, in order to optimize resource handling.

Example 2. Consider a program consisting of a main fragment in *HL* invoking the **external** *SL* function applyToInt with the *HL* function add1 as an argument; the call is meaningful because add1 does not invoke **callcc**. Only the *SL* type of the external function is given to the *HL* program which is separately compiled without access to the detailed effect annotations inferred from the code of the *SL* fragment.

SL fragment apply Tolnt:

$$\lambda f: Int \rightarrow Int. succ (f 42)$$

The result of its separate compilation into *RL*, which uses stack allocation (for details of the translation we refer the reader to Section 5) is

applyToInt = 
$$\Lambda t \le S$$
.  $\lambda^S f$ :Int  $\frac{S}{t}$  Int. let  $S = 0$  f 42 in  $S = 0$  succ  $S = 0$  suc

HL fragment main:

**external**(SL) applyToInt : (Int  $\rightarrow$  Int)  $\rightarrow$  Int in let add1 =  $\lambda x$  : Int. succ x in applyToInt add1

The result of its separate compilation into RL is

$$\begin{split} \text{main} &= \lambda^{\text{H}} \, \text{applyToInt:} \, \forall t \leq S. \, \left( \text{Int} \, \frac{s}{t} \, \text{Int} \right) \, \frac{s}{\emptyset} \, \text{Int.} \\ \textbf{let}^{\text{H}} \\ & \text{applyToInt\_H} = \langle \Delta t \leq S. \\ & \lambda^{\text{H}} \, \text{f:Int} \, \frac{H}{t} \, \text{Int.} \\ & \textbf{let}^{\text{H}} \, \text{f\_S} = \langle \lambda^{S} \, x \colon \text{Int.} \, \textbf{use}^{\text{H}} \, \left( \mathfrak{G} \, f \, x \right) \rangle^{\text{H}} \\ & \text{in} \, \, \textbf{use}^{S} \, \left( \mathfrak{G} \, \left( \text{applyToInt[t]} \right) \, f\_S \right) \rangle^{\text{H}} \\ & : \forall t \leq S. \, \left( \text{Int} \, \frac{H}{t} \, \text{Int} \right) \, \frac{H}{\emptyset} \, \text{Int} \\ & \text{add1} \qquad = \langle \lambda^{H} \, x \colon \text{Int.} \, \mathfrak{G} \, \text{succ} \, x \rangle^{\text{H}} \\ & : \text{Int} \, \frac{H}{\emptyset} \, \text{Int} \\ & \text{in} \, \, \mathfrak{G} \, \text{applyToInt\_H}[\emptyset] \, \text{add1} \\ & : \left( \forall t \leq S. \, \left( \text{Int} \, \frac{S}{t} \, \text{Int} \right) \, \frac{S}{\emptyset} \, \text{Int} \right) \, \frac{H}{\emptyset} \, \text{Int} \end{split}$$

The translation infers polymorphic effect types using a simplified version<sup>2</sup> of standard effect inference [23]. The resource annotations are fixed by the source language; the type of an external SL function in an HL program is annotated with the SL resources. In the code produced after translation the external functions are coerced to match the resources of HL using automatically generated wrappers. In the above code, the parameter f of applyTolnt\_H is wrapped to f\_S before passing it to applyTolnt; the function of the wrapper is to switch from the stack allocation discipline used by SL to heap allocation before invoking the code for f, and resume stack allocation upon return. Dually, the call to applyTolnt itself is wrapped to enable stack allocation inside HL code.

<sup>&</sup>lt;sup>2</sup> As presented here our system does not keep track of regions associated with effects.

Since the full RL type of the SL fragment is not available to it, the effect inference must conseratively approximate the effects of the SL functions. It treats the external applyToInt in the HL fragment as an effect-polymorphic parameter in order to allow its invocations with arguments with different effects. The price we pay for inference with this polymorphism in the case of separate compilation is that we assume that the effects of these invocations are the maximal allowed with the resources shared between the languages (in Example 2 we lose no precision since SL has no effects, but the approximation is reflected in the effect annotation  $\emptyset$  of the type of the parameter of main). The following code, constructed mechanically given the inferred and expected types of applyToInt, coerces the actual type of applyToInt to the approximation used in the typing of main and performs the top-level application, thus linking the modules.

$$\begin{split} \textbf{let}^{H} \\ & \text{applyToInt\_Glue} = \langle \Lambda t \leq S. \ \lambda^{S} \ f \colon \text{Int.} \ \stackrel{S}{\underset{\theta}{\longrightarrow}} \ \text{Int.} \ \langle \text{@ applyToInt}[t] \ f \rangle_{\emptyset} \rangle^{H} \\ & : \forall t \leq S. \ (\text{Int} \ \frac{S}{\underset{\theta}{\longrightarrow}} \ \text{Int}) \ \frac{S}{\underset{\theta}{\longrightarrow}} \ \text{Int} \end{split}$$

in @ main applyToInt\_Glue

More precise inference of the resulting effects is possible when the external function is a pre-compiled library routine whose *RL* type (with its precise effect annotations) is available when compiling main. In those cases we can take advantage of the letpolymorphism in inferring a type of main (in a setting similar to that of Example 1). However even the approximated effects obtained during separate compilation carry information that can be exploited for the optimization of inter-language calls, observing that the range of effects of a function is limited by the resources of its source language. In Example 2, after inlining and applying results of Section 4.4 (Theorem 2), the code for main can be optimized to eliminate the unnecessary switch to heap allocation in the instance of f\_S. This yields

$$\begin{split} \mathsf{main} &= \langle \lambda^\mathsf{H} \ \mathsf{applyToInt} \colon \forall \mathsf{t} \leq \mathsf{S}. \ (\mathsf{Int} \ \frac{\mathsf{S}}{\mathsf{t}} \ \mathsf{Int}) \ \frac{\mathsf{S}}{\emptyset} \ \mathsf{Int}. \\ & \mathsf{let}^\mathsf{H} \\ & \mathsf{add1} \ = \langle \lambda^\mathsf{H} \ \mathsf{x} \colon \mathsf{Int}. \ \mathsf{Q} \ \mathsf{succ} \ \mathsf{x} \rangle^\mathsf{H} \\ & \mathsf{add1} \mathcal{S} &= \langle \lambda^\mathsf{S} \ \mathsf{x} \colon \mathsf{Int}. \ \mathsf{Q} \ \mathsf{succ} \ \mathsf{x} \rangle^\mathsf{H} \\ & \mathsf{in} \ \mathsf{use}^\mathsf{S} \ (\mathsf{Q} \ (\mathsf{applyToInt}[\emptyset]) \ \mathsf{add1} \mathcal{S}) \rangle^\mathsf{H} \end{split}$$

Thus the HL function add1 has been effectively specialized for the stack allocation strategy used by SL.

Example 3. Another optimization is merging of regions with the same resource requirements, illustrated on the following HL code fragment.

external(
$$SL$$
) intFn : Int  $\rightarrow$  Int in intFn (intFn 42)

which is naively translated to the RL function (shown after inlining of the parameter wrapper)

$$\begin{split} \Lambda t \! \leq \! S. \, \lambda^{H} \, \text{intFn:Int} \, & \xrightarrow{S} \, \text{Int.} \\ & \text{let}^{H} \, x = \langle \text{use}^{S} \, (\text{@intFn 42}) \rangle^{H} \\ & \text{in use}^{S} \, (\text{@intFn x}) \end{split}$$

After combining the two  $use^{S}$  (·) constructs the equivalent RL term is

$$\begin{split} \Lambda t &\leq S. \; \lambda^{H} \, \text{intFn:Int} \, \frac{s}{t} \, \text{Int.} \\ & \quad \text{use}^{S} \, \left( \, \text{let}^{S} \, \, x = \langle \text{@intFn 42} \rangle^{S} \, \, \text{in @intFn x} \right) \end{split}$$

A generalization of this transformation makes possible lifting of  $\mathbf{use}^r$  (·) constructs out of a loop when the resources r are sufficient for all effects of the loop. Since in general a resource wrapper must restore resources upon return, a tail call moved into its scope effectively becomes non-tail; thus lifting a wrapper's scope over a recursive tail call is only useful when the wrapper is lifted out of the enclosing function as well, i.e. out of the loop.

# 4 Semantics of RL

#### 4.1 Static Semantics

Correctness of resource use is ensured by the type system shown in Figure 3, which keeps track of the resources necessary for the evaluation of a term and a conservative estimate of the effects of the evaluation.

An effect environment  $\Delta$  specifies the resource bounds of effect variables introduced by effect abstractions and effect-polymorphic types. The rules for effect sequents reflect the dependence of effects on resources (in this language this boils down to the dependence of the call/cc effect CC on the heap allocation resource H) and form the basis of effect polymorphism. The function MaxEff yields the maximal effect possible with a given resource; in this system we have  $MaxEff(S) = \emptyset$  and MaxEff(H) = CC. Rule (Eff-max) effectively states that the resource r' can be used instead of resource r if r' provides for all effects possible under r.

In the sequents assigning types to values and terms the type environment  $\Gamma$  maps free variables to types. Type judgments for values associate with a value v and a pair of environments  $\Delta$  and  $\Gamma$  only a type  $\sigma$ , since values have no effects and therefore their evaluation requires no resources of the kind we control. The function  $\theta$  maps constants to their predefined types.

Sequents for terms have the form  $r; \Delta; \Gamma \vdash_e e : \frac{1}{\mu}\sigma$ , where r represents the available allocation resource,  $\sigma$  is the type of e, and  $\mu$  represents the effects of its evaluation. Rules (Exp-let) and (Exp-val) establish the correspondence between the resource annotations in these constructs and the currently available allocation resource; the effect of lifting a value to a term is none, while the effect of sequencing two computations via let is the union of their effects. Any effect allowed with the current resource may be added to the effects of a term using rule (Exp-spurious).

The central novelty is the  $\mathbf{use}^{r'}(\cdot)$  construct for resource manipulation; its typing rule (Exp-use) imposes the crucial restriction that the effect  $\mu$  of the term e must be

# EFFECT ENVIRONMENT FORMATION (Env. off ampty) (Env. off ampty)

# TYPE ENVIRONMENT FORMATION

#### **EFFECTS**

$$\begin{array}{c} \textbf{(Eff-empty)} \\ \frac{\vdash^{\Delta} \Delta}{\Delta \vdash^{\mu} \emptyset \leq r} \end{array} \qquad \begin{array}{c} \textbf{(Eff-CC)} \\ \stackrel{\vdash^{\Delta} \Delta}{\Delta \vdash^{\mu} \mathsf{CC} \leq \mathsf{H}} \end{array} \\ \\ \textbf{(Eff-var)} \\ \frac{\vdash^{\Delta} \Delta \Delta(t) = r}{\Delta \vdash^{\mu} t \leq r} \qquad \begin{array}{c} \textbf{(Eff-combine)} \\ \frac{\vdash^{\mu} \mu' \leq r, \ \mu'' \leq r}{\Delta \vdash^{\mu} \mu' \lor \mu'' \leq r} \end{array} \\ \\ \textbf{(Eff-max)} \\ \frac{\Delta \vdash^{\mu} \mu \leq r \quad \Delta \vdash^{\mu} \mathit{MaxEff}(r) \leq r'}{\Delta \vdash^{\mu} \mu \leq r'} \end{array}$$

#### VALUES

$$\begin{array}{c} \textbf{(Val-const)} \\ \underline{\Delta \vdash^{\Gamma} \Gamma} \\ \underline{\Delta ; \Gamma \vdash^{\nu} c : \theta(c)} \end{array} \qquad \begin{array}{c} \textbf{(Val-var)} \\ \underline{\Delta \vdash^{\Gamma} \Gamma \Gamma(x) = \sigma} \\ \underline{\Delta \vdash^{\Gamma} \Gamma \Gamma(x) = \sigma} \\ \\ \underline{A \vdash^{\Gamma} \Gamma \Delta \vdash^{\sigma} \sigma} \\ \underline{r ; \Delta ; \Gamma_{x}, \ x : \sigma \vdash^{e} e : \underline{\mu} \sigma'} \\ \underline{\Delta ; \Gamma_{x} \vdash^{\nu} \lambda^{r} \ x : \sigma \cdot e : \sigma \stackrel{r}{\xrightarrow{\rho} \sigma'} } \\ \underline{A ; \Gamma_{x} \vdash^{\nu} \lambda^{r} \ x : \sigma \cdot e : \sigma \stackrel{r}{\xrightarrow{\rho} \sigma'} } \\ \underline{A \vdash^{\Gamma} \Gamma \Delta_{t}, \ t \leq r ; \Gamma \vdash^{\nu} v : \sigma} \\ \underline{\Delta ; \Gamma \vdash^{\nu} \Lambda t \leq r \cdot v : \forall t \leq r \cdot \sigma} \\ \underline{A \vdash^{\Gamma} \Gamma \Delta_{t}, \ t \leq r \cdot \sigma \quad \Delta \vdash^{\mu} \mu \leq r} \\ \underline{A \vdash^{\Gamma} \Gamma \Delta_{t}, \ r \vdash^{\nu} x [\mu] : [\mu/t] \sigma} \end{array}$$

$$(\textbf{Typ-basic}) \qquad \begin{matrix} (\textbf{Typ-fun}) \\ (\textbf{Zp-fun}) \\ \underline{\Delta \vdash^{\mu} \Delta} \end{matrix} \qquad \begin{matrix} \underline{\Delta \vdash^{\mu} \mu \leq r \quad \Delta \vdash^{\sigma} \sigma, \ \sigma'} \\ \underline{\Delta \vdash^{\sigma} \sigma \stackrel{\tau}{\rightarrow} \sigma'} \end{matrix}$$

$$(\textbf{Typ-cont}) \qquad \qquad \begin{matrix} \underline{\Delta \vdash^{\sigma} \sigma \quad \emptyset \vdash^{\mu} CC \leq r} \\ \underline{\Delta \vdash^{\sigma} \sigma \quad^{\tau} cont} \end{matrix}$$

$$(\textbf{Typ-poly}) \qquad \qquad \begin{matrix} \underline{\vdash^{\Delta} \Delta \quad \Delta_{t}, \ t \leq r \vdash^{\sigma} \sigma} \\ \underline{\Delta \vdash^{\sigma} \forall t \leq r. \ \sigma} \end{matrix}$$

#### **TERMS**

$$\begin{split} &(\textbf{Exp-let}) \\ & r; \Delta; \Gamma \vdash^e e : \frac{-\sigma}{\mu} \quad r; \Delta; \Gamma_x, \ x : \sigma \vdash^e e' : \frac{-\sigma'}{\mu'} \sigma' \\ & r; \Delta; \Gamma \vdash^e \textbf{let}^r \ x = e \ \textbf{in} \ e' : \frac{-\sigma'}{\mu' \vee \mu'} \sigma' \\ & \frac{(\textbf{Exp-val})}{r; \Delta; \Gamma \vdash^e \langle v \rangle^r : \frac{-\sigma}{\emptyset}} \\ & \frac{\Delta; \Gamma \vdash^e v : \sigma}{r; \Delta; \Gamma \vdash^e e : -\sigma \quad \Delta \vdash^\mu \mu' \leq r} \\ & \frac{r; \Delta; \Gamma \vdash^e e : -\sigma \quad \Delta \vdash^\mu \mu' \leq r}{r; \Delta; \Gamma \vdash^e e : -\sigma \quad \Delta \vdash^\mu \mu \leq r} \\ & \frac{(\textbf{Exp-use})}{r; \Delta; \Gamma \vdash^e \textbf{use}^r (e) : \frac{-\sigma}{\mu'}} \\ & \frac{(\textbf{Exp-use})}{r; \Delta; \Gamma \vdash^e \textbf{use}^r (e) : -\sigma} \\ & \frac{(\textbf{Exp-app})}{r; \Delta; \Gamma \vdash^e \textbf{Q} \ x \ x' : -\sigma} \\ & \frac{(\textbf{Exp-callcc})}{r; \Delta; \Gamma \vdash^e \textbf{Callcc} \ x : \frac{-\sigma}{\mu'} \sigma} \\ & \frac{(\textbf{Exp-callcc})}{r; \Delta; \Gamma \vdash^e \textbf{callcc} \ x : \frac{-\sigma}{\mu'} \sigma} \end{split}$$

$$\frac{\Delta \vdash^r \Gamma \quad \Delta \vdash^\sigma \sigma' \quad \Gamma(x) = \sigma \quad r \mathsf{cont} \quad \Gamma(x') = \sigma}{r; \Delta; \Gamma \vdash^e \mathsf{throw}[\sigma'] \ x \ x' : \frac{}{MaxEff(r)} \sigma'}$$

**Fig. 3.** The *RL* type system

supported by the resource r available before the alternative resource r' is selected. This ensures the correctness of the propagation of  $\mu$  outside the scope of the  $\mathbf{use}^{r'}$  ( $\cdot$ ).

The rules for application and **callcc** set the correspondence between the available resource and the resource required by the invoked function. In addition, (Exp-callcc) and (Exp-throw) specify that the continuation type is annotated with the same resource, which is needed by the context captured in the continuation and therefore must be matched when it is reactivated. The effect of evaluating a **callcc** includes CC, while the effect of a **throw** is that of the rest of the computation, which we estimate as the maximal possible with the current resource.

By induction on the structure of a typing derivation it follows that if a term has a type in a given environment, it has exactly one type, and the presence of type annotations allows its effective computation, *i.e.* there exists a function EffTypeOf such that

$$\textit{EffTypeOf}\ (r,\ \Delta,\ \Gamma,\ e) = \langle \mu,\sigma\rangle \ \text{if and only if}\ r;\Delta;\Gamma\ \vdash^e e: -\sigma.$$

We will also use the function TypeOf with the same arguments, returning the type  $\sigma$  only.

# 4.2 Dynamic Semantics

The operational semantics of RL (Figure 4) is defined by means of a variant of the tail-call-safe  $C_aEK$  machine (Flanagan  $et\ al.$  [4]). The machine configuration is a tuple  $\langle e, E, O, \rho \rangle$  where e is the current term to be evaluated, E is the environment mapping variables to machine values, E is a heap of objects (closures), and E is a tuple of machine resources. Depending on the allocation strategy used, E is either a continuation stack E, recording (as in the original E is a current continuation E is a continuation heap E. In the latter form E is a continuation handle and E is a mapping from E to activation records which offers non-sequential access. In neither case does a function application (app) perform additional allocations of activation records, so both strategies are tail-call safe.

Machine values are either small constants or pointers into other structures where larger objects are allocated. All closures are allocated on the heap (the function  $\gamma$  at the bottom of the figure shows the details).

The activation records created when evaluating a  $\mathbf{let}^{r}$ -expression may be allocated either on the continuation heap K (transition rule  $(\mathbf{let}^{H})$ ) or on the continuation stack S (rule  $(\mathbf{let}^{S})$ ). An activation record represents a continuation, and in our small language there are only three possibilities: the computation either halts or continues by binding a variable to a computed value or by restoring a resource. Rules  $(\mathbf{val}^{H})$  and  $(\mathbf{val}^{S})$  perform the binding, depending on the allocation mode.

The evaluation of  $\mathbf{use}^r$  (e) selects the activation record allocation strategy for e, e.g.  $\mathbf{use}^s$  (e) selects stack-based allocation for e (transition rule ( $\mathbf{use}^s$ )). When the current allocation resource is already r we define  $\mathbf{use}^r$  (·) as a no-op; if a change of resource is performed, an activation record is pushed on (the top of) the new allocation resource. Correspondingly, heap-based allocation is restored by transition rule ( $\mathsf{resume}^\mathsf{H}$ ) after the evaluation of e.

#### SEMANTIC DOMAINS $Machine Val \ni w ::= \mathsf{Const}\ c \mid \mathsf{Ptr}\ h \mid \mathsf{Cont}\ k$ machine values $E \in Var \rightarrow Machine Val$ environment $h \in HeapLocs$ heap locations $Object ightarrow o ::= Closure \langle x, e, E \rangle \mid TyAbs \langle t, r, v \rangle$ closures (objects) $O \in HeapLocs \rightarrow Object$ object heap $k \in ContHandles$ continuation handles $ActRcd \ni a ::=$ Bind $\langle x, e, E, k \rangle \ | \$ Resume $S \ | \$ Halt activation records $K \in ContHandles \rightarrow ActRcd$ activation record heap $S ::= \ \mathsf{Bind} \ \langle x, e, E, S \rangle \ | \ \mathsf{Resume} \ \langle k, K \rangle \ | \ \mathsf{Halt}$ activation record stack TRANSITION RULES $\langle \mathbf{0} x_1 x_2, E, O, \rho \rangle \mapsto_1 \langle e', E'[x' \mapsto E(x_2)], O, \rho \rangle$ (app) where $E(x_1) = \operatorname{Ptr} h$ , $O(h) = \operatorname{Closure} \langle x', e', E' \rangle$ FOR HEAP-ALLOCATED ACTIVATION RECORDS (let<sup>H</sup>) $\langle \mathsf{let}^\mathsf{H} \ x = e_1 \ \mathsf{in} \ e_2, \ E, \ H, \ \langle k, \ K \rangle \rangle \ \mapsto_1$ $\langle e_1, E, H, \langle k', K[k' \mapsto \mathsf{Bind} \langle x, e_2, E|_{FV(e_2)-x}, k \rangle] \rangle \rangle$

$$\begin{aligned} \text{(val}^{\mathsf{H}}) & & & & & & & & & & & & & & & & \\ & & & & & & & & & & & & \\ & & & & & & & & & & \\ & & & & & & & & & \\ & & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ & & & \\ & & \\ & & & \\ & &$$

$$\begin{array}{ll} \text{(throw)} & & \langle \mathsf{throw}[\sigma] \; x_1 \; x_2, \; E, \; H, \; \langle k, \; K \rangle \rangle \mapsto_1 \langle e', \; E'[x' \mapsto E(x_2)], \; O, \; \langle k', \; K \rangle \rangle \\ & & \text{where} \; E(x_1) = \mathsf{Cont} \; k_1, \; K(k_1) = \mathsf{Bind} \; \langle x', e', E', k' \rangle \end{array}$$

$$(\mathsf{use}^\mathsf{S}) \qquad \qquad \langle \mathsf{use}^\mathsf{S} \ (e), \, E, \, H, \, \langle k, \, K \rangle \rangle \mapsto_1 \langle e, \, E, \, H, \, \langle \mathsf{Resume} \ \langle k, K \rangle \rangle \rangle$$

### FOR STACK-ALLOCATED ACTIVATION RECORDS

$$\langle \mathsf{let}^{\mathsf{S}} \rangle \qquad \langle \mathsf{let}^{\mathsf{S}} \ x = e_1 \ \mathsf{in} \ e_2, \ E, \ H, \ \langle S \rangle \rangle \ \mapsto_1 \ \langle e_1, \ E, \ H, \ \langle \mathsf{Bind} \ \langle x, e_2, E|_{FV(e_2)-x}, S \rangle \rangle \rangle$$

$$(\mathsf{use}^\mathsf{H}) \qquad \qquad \langle \mathsf{use}^\mathsf{H} \left( e \right), \, E, \, H, \, \langle S \rangle \rangle \mapsto_1 \langle e, \, E, \, H, \, \langle k, \, [k \mapsto \mathsf{Resume} \, S] \rangle \rangle$$

$$(\mathsf{resume}^{\mathsf{H}}) \qquad \langle \langle v \rangle^{\mathsf{S}}, \, E, \, H, \, \langle \mathsf{Resume} \, \langle k, K \rangle \rangle \rangle \mapsto_1 \langle \langle v \rangle^{\mathsf{H}}, \, E, \, H, \, \langle k, \, K \rangle \rangle$$

#### REPRESENTATION OF VALUES

(val<sup>H</sup>)

Fig. 4. Semantics of RL

Another no-op is the increase of effect sets  $\langle \cdot \rangle_{\mu}$  which only serves type-checking purposes.

#### 4.3 Soundness of the Type System

The type system maintains the property that the effects of well-typed programs are possible with their available resources, formalized in the following statement, proved by induction on the typing derivation.

**Lemma 1.** If r;  $\Delta$ ;  $\Gamma \vdash^e e : \frac{1}{\mu}\sigma$  is a valid typing judgment, then  $\Delta \vdash^{\mu} \mu \leq r$ .

Semantically this behavior of well-typed programs is expressed as soundness with respect to resource use, extending the standard soundness for safety of the type system, in the following theorem.

**Theorem 1.** If  $r; \emptyset; \emptyset \vdash e : \frac{}{\mu}\sigma$ , then the configuration  $\langle e, \emptyset, \emptyset, \mathsf{Halt}^r \rangle$  either diverges or evaluates to the configuration  $\langle \langle v \rangle^r, E, O, \langle \mathsf{Halt}^r \rangle \rangle$  (for some v, E and O), where  $\mathsf{Halt}^\mathsf{S} \stackrel{\triangle}{=} \langle \mathsf{Halt} \rangle$ , and  $\mathsf{Halt}^\mathsf{H} \stackrel{\triangle}{=} \langle k, K \rangle$  for some k and K such that  $K(k) = \mathsf{Halt}$ .

This result is a corollary of the standard properties of progress and subject reduction of the system, the proofs of which we sketch below. To simplify the proofs, we introduce a type-annotated version of the semantics, which maintains type information embedded in the runtime representation. Thus the representation of an abstraction in the type-annotated version is

$$\gamma(\lambda^r x : \sigma. e, E, O) = \langle \mathsf{Ptr} h, O[h \mapsto \mathsf{Closure}' \langle r, x, \sigma, e, E|_{FV(e)-x} \rangle] \rangle$$

In addition, the runtime environment E is extended to keep the type of each value in its codomain; the value component of E is denoted by  ${}^{V}E$  and the type component by  ${}^{T}E$ . The following definitions are helpful in defining typability of configurations.

**Definition 1.** The bottom  $bot(\rho)$  of an allocation resource  $\rho$  is defined as follows:

- 1. if  $\rho = \langle S \rangle$ , then  $bot(\rho) = bot(S')$ , if  $S = Bind \langle x', e', E', S' \rangle$ , and  $bot(\rho) = S$  otherwise;
- 2. if  $\rho = \langle k, K \rangle$ , then  $bot(\rho) = bot(\langle k', K \rangle)$ , if  $K(k) = Bind \langle x', e', E', k' \rangle$ , and  $bot(\rho) = K(k)$  otherwise.

**Definition 2.** The outermost continuation heap  $outerCont(\rho)$  reachable from allocation resource  $\rho$  is

- 1.  $K \text{ if } \rho = \langle k, K \rangle \text{ and } bot(\rho) = \mathsf{Halt};$
- 2.  $outerCont(\langle S \rangle)$  if  $\rho = \langle k, K \rangle$  and  $bot(\rho) = Resume S$ ;
- 3.  $\emptyset$ , if  $\rho = \langle S \rangle$  and  $bot(\rho) = Halt$ ;
- 4.  $outerCont(\langle k, K \rangle)$  if  $\rho = \langle S \rangle$  and  $bot(\rho) = Resume \langle k, K \rangle$ .

**Definition 3.** A configuration closed in type environment  $\Gamma$  is typable under resource r with a result type  $\sigma$  and an effect  $\mu$ , written r;  $\Gamma \vdash \langle e, E, O, \rho \rangle : \frac{1}{\mu}\sigma$ , if for some  $\sigma'$ ,  $\mu'$ 

```
1. Dom(\Gamma) \cap Dom(E) = \emptyset; and

2. r; \emptyset; \Gamma, {}^{T}E \vdash e : \frac{r}{\mu^{T}}\sigma'; and

3. \Gamma \vdash P(\rho, E, O) \in \sigma' \xrightarrow{r} \sigma; and

4. for each \ x \in Dom(E),

(a) if \ {}^{V}E(x) = Const \ c, then \ {}^{T}E(x) = \theta(c);

(b) if \ {}^{V}E(x) = Ptr \ h \ and \ O(h) = Closure' \ \langle r_1, x_1, \sigma_1, e_1, E_1 \rangle, then \ \emptyset; {}^{T}E_1 \vdash P \lambda^{r_1} x_1 : \sigma_1. e_1 : {}^{T}E(x), and similarly for type abstractions;

(c) if \ {}^{V}E(x) = Cont \ k, then \ {}^{T}E(x) = \sigma_1 \ {}^{r_1}cont \ and \ \Gamma \vdash P \langle k, outerCont(\rho) \rangle, E, O \in \sigma_1 \xrightarrow{r_1} \sigma'_1

and \ \mu = \mu_1 \lor \mu'_1, for some \ \sigma'_1 \ and \ \mu'_1,

and \ \Gamma \vdash P \langle \rho, E, O \rangle \in \sigma' \xrightarrow{r} \sigma \ if
```

1. r = S and  $\rho = \langle \mathsf{Halt} \rangle$  (i.e. an empty stack) and  $\sigma = \sigma'$  and  $\mu = \emptyset$ ; or
2. r = S and  $\rho = \langle \mathsf{Bind} \langle x_1, e_1, E_1, S_1 \rangle \rangle$  and S;  $\Gamma$ ,  $x_1 : \sigma' \vdash^e \langle e_1, E_1, O, S_1 \rangle : \frac{1}{\mu}\sigma$ ; or
3. r = S and  $\rho = \langle \mathsf{Posumo} \rangle \langle h', K' \rangle \rangle$  and  $\Gamma \vdash^e \langle h', K' \rangle \rangle F$ .  $O \setminus C$   $\sigma' \stackrel{\mathsf{I}}{\hookrightarrow} \rangle \sigma$ 

3.  $r = \mathsf{S} \ and \ \rho = \langle \mathsf{Resume} \ \langle k', K' \rangle \rangle \ and \ \Gamma \ \vdash^\rho \ \langle \langle k', K' \rangle, E, O \rangle \in \sigma' \overset{\mathsf{H}}{\underset{\mu}{\hookrightarrow}} \sigma,$ 

and similarly for r = H.

Note that the environment may contain reachable variables bound to continuations even when the current allocation resource is a stack. Type correctness of these continuations cannot be verified with the stack resource, instead we have to find the corresponding continuation heap. However in this case the type system guarantees that the only continuation heap to which there are references in the environment is the outermost continuation heap, if such exists. The reason is that although it is possible to switch to heap allocation after executing in stack allocation mode, there are no invocations of **callcc** allowed since they would introduce the CC effect, which is not possible under the stack resource (cf. typing rule (Exp-use) in Figure 3).

We can now formulate the progress and subject reduction properties.

**Lemma 2** (**Progress**). If r;  $\emptyset \vdash^c \langle e, E, O, \rho \rangle : \frac{1}{\mu}\sigma$  where r corresponds to  $\rho$  (i.e. r = S if  $\rho = \langle S \rangle$ , r = H if  $\rho = \langle k, K \rangle$ ), and  $\rho \neq H$ alt r, then there exists C such that  $\langle e, E, O, \rho \rangle \mapsto_1 C$ .

**Lemma 3** (Subject reduction). If  $C = \langle e, E, O, \rho \rangle$  and  $r; \emptyset \vdash^c C : \frac{}{\mu}\sigma$  where r corresponds to  $\rho$ , and  $C \mapsto_1 C' = \langle e', E', O', \rho' \rangle$ , then  $r'; \emptyset \vdash^c C' : \frac{}{\mu'}\sigma$  where r' corresponds to  $\rho'$ ,  $\mu = \mu' \vee \mu'_1$ , and the rule for this transition is (callcc) only if  $\mu = \mathsf{CC} \vee \mu''$ , for some  $\mu'_1$  and  $\mu''$ .

In brief, in the case when  $e \neq \langle v \rangle^r$ , the proofs proceed by examining the structure of the typing derivation for  $r; \emptyset; \Gamma, \ ^T\!E \vdash^e e : \frac{}{\mu^r}\sigma';$  together with condition 4 of Definition 3 this yields that the values in the environment and on the heaps have the correct shape for the appropriate transition rule. In the case when e has the form  $\langle v \rangle^r$  the proofs inspect the structure of the derivation of  $\Gamma \vdash^e \langle \rho, E, O \rangle \in \sigma' \stackrel{r}{\hookrightarrow} \sigma$ , which parallels the decision tree for the transition rules (val) and (resume) and the halting state.

#### 4.4 Resource Transformations

Effect inference and type correctness with respect to resource use allow the compiler to modify the continuation allocation strategy of a program fragment and preserve its meaning. The following definitions adapt the standard notions of ordering and observational equivalence of open terms to the resource-based system.

**Definition 4.** A context C is a term with a hole  $\bullet$ ; the result of placing a term e in the hole of C is denoted by C[e] and may result in capturing effect and lambda variables free in e. The hole of a context C is of type  $(r, \Delta, \Gamma) \Rightarrow \frac{1}{\mu}\sigma$  if C[e] is typeable whenever  $r; \Delta; \Gamma \vdash e : \frac{1}{\mu}\sigma$ .

**Definition 5.** S;  $\Delta$ ;  $\Gamma \vdash e \sqsubseteq e' : \frac{}{\mu}\sigma$  if for all contexts C with hole of type  $(r, \Delta, \Gamma) \Rightarrow \frac{}{\mu}\sigma$ , all typed environments E closing C[e] and heaps O closing E, and continuation stacks S, the configuration  $\langle C[e'], E, O, \langle S \rangle \rangle$  converges if  $\langle C[e], E, O, \langle S \rangle \rangle$  converges. Furthermore, S;  $\Delta$ ;  $\Gamma \vdash e e \approx e' : \frac{}{\mu}\sigma$  if S;  $\Delta$ ;  $\Gamma \vdash e e \sqsubseteq e' : \frac{}{\mu}\sigma$  and S;  $\Delta$ ;  $\Gamma \vdash e e' \sqsubseteq e : \frac{}{\mu}\sigma$ .

One possible optimization is the conversion of heap-allocating code to stack-based strategy provided the code does not invoke **callcc** or **throw**, as per the following theorem.

**Theorem 2.** If H;  $\Delta$ ;  $\Gamma \vdash^{e} e : \overline{_{\emptyset}}\sigma$ , then S;  $\Delta$ ;  $\Gamma \vdash^{e} \mathbf{use}^{\mathsf{H}}(e) \approx StkCont_{\Delta}(e; \Gamma) : \overline{_{\emptyset}}\sigma$ , where StkCont is the transformation defined as follows.

```
StkCont_{\Delta}\left(\left\langle v\right\rangle^{\mathsf{H}};\;\Gamma\right)=\left\langle v\right\rangle^{\mathsf{S}} StkCont_{\Delta}\left(\left\langle e\right\rangle_{\mu};\;\Gamma\right)=\left\langle StkCont_{\Delta}\left(e;\;\Gamma\right)\right\rangle_{\mu} StkCont_{\Delta}\left(\mathsf{use}^{\mathsf{H}}\left(e\right);\;\Gamma\right)=StkCont_{\Delta}\left(e;\;\Gamma\right) StkCont_{\Delta}\left(\mathsf{use}^{\mathsf{S}}\left(e\right);\;\Gamma\right)=e StkCont_{\Delta}\left(\mathsf{0}\;x_{1}\;x_{2};\;\Gamma\right)=\mathsf{let}^{\mathsf{S}}\;x_{1}'=\left\langle \lambda^{\mathsf{S}}\;x_{2}':\Gamma(x_{2}).\;\mathsf{use}^{\mathsf{H}}\left(\mathsf{0}\;x_{1}\;x_{2}'\right)\right\rangle^{\mathsf{S}} \mathsf{in}\;\mathsf{0}\;x_{1}'\;x_{2} StkCont_{\Delta}\left(\mathsf{let}^{\mathsf{H}}\;x=e_{1}\;\mathsf{in}\;e_{2};\;\Gamma\right)=\mathsf{let}^{\mathsf{S}}\;x=StkCont_{\Delta}\left(e_{1};\;\Gamma\right) \mathsf{in}\;StkCont_{\Delta}\left(e_{2};\;\Gamma_{x},x:TypeOf\left(\mathsf{H},\Delta,\Gamma,e_{2}\right)\right)
```

#### 5 Translation from HL to RL

Programs in language  $\mathcal{L} \in \{HL, SL\}$  are translated into RL by an algorithm shown in Figure 5. The algorithm infers the effect and resource annotations of a term using fairly standard techniques. It is presented in the form of an inference system for judgments of the form  $\Delta$ ;  $\Gamma \vdash_{\mathcal{L}} e_{HL} \Rightarrow \Delta' \vdash_{e} : \frac{1}{\mu} \sigma$ , where  $e_{HL}$ ,  $\Delta$ , and  $\Gamma$  are inputs corresponding respectively to the  $\mathcal{L}$  term to translate (also overloaded to HL top-level programs) and the inherited effect and type environments, initially empty. The outputs of the translation are e,  $\Delta'$ ,  $\mu$ , and  $\sigma$ , which stand for the translated term, the inferred effect environment, and the effect and type of e in environments  $\Delta'$  and  $\Gamma$ ; thus the output of the algorithm satisfies H;  $\Delta'$ ;  $\Gamma \vdash_{e} e : \frac{1}{\mu} \sigma$ . The function  $\mathcal{R}$  maps a language name to the resources available to a program in this language:  $\mathcal{R}(HL) = H$  and  $\mathcal{R}(SL) = S$ .

$$\begin{aligned} & (\text{Translate-external}) \\ & \sigma' = CloseAll \left( \text{Max}^{\$}(\text{Annotate}^{\$}(\tau, Dom\left(\Delta\right))), \$ \right) \\ & \Delta; \Gamma \vdash_{\text{th.}} \text{external}(\text{SL}) \ x : \tau \text{ in p} \\ & \Rightarrow \Delta' \vdash \lambda^{\text{H}} \ x : \sigma', \text{ let}^{\text{H}} \ x = Wrap_{\$}^{\text{H}}(\bullet, x, \sigma') \text{ in } e' : \frac{1}{\$}(\sigma') \xrightarrow{\frac{1}{\mu}} \sigma) \end{aligned}$$
 where 
$$& \text{Annotate}^{\text{T}}(\beta, V) = \beta \\ & \text{Annotate}^{\text{T}}(\tau, V) = \delta \\ & \text{Annotate}^$$

Fig. 5. Typed translation from HL to RL

Several auxiliary functions are shown in the figure, and the definitions of several simpler functions are as follows. The lub of two resources is defined by  $r \sqcup r = r$  and  $S \sqcup H = H$ . The function  $\sqcap$  for merging two effect environments is defined as  $(\Delta_1 \sqcap \Delta_2)(t) = \Delta_1(t) \sqcup \Delta_2(t)$  if  $t \in Dom(\Delta_1) \cap Dom(\Delta_2)$ , and  $(\Delta_1 \sqcap \Delta_2)(t) = \Delta_i(t)$  on the rest of  $Dom(\Delta_1) \cup Dom(\Delta_2)$ . The free effect variables of a type  $\sigma$  are denoted by  $fev(\sigma)$ ; the function  $Close(\sigma, \Delta, \Gamma)$  returns the pair  $\langle \forall t_i \leq \Delta(t_i), \sigma, \Delta \setminus_{\{\overline{t_i}\}} \rangle$ , where  $\{\overline{t_i}\} = fev(\sigma) - fev(\Gamma)$ , and similarly we have  $CloseAll(\sigma, r) = \forall \overline{t_i} \leq r, \sigma$  where  $\{\overline{t_i}\} = fev(\sigma)$ .

Separately compiled **external** functions are treated as parameters of the compiled HL fragment and are wrapped to convert the HL resources (continuation heap) to SL resources (continuation stack). The wrapping is performed by an auxiliary function invoked as  $Wrap_r^{r'}(C, x, \sigma)$ , which produces a term coercing x from type  $\sigma$  to type  $ConvertType_r^{r'}(\sigma)$  with resource annotations r' in place of r, and places it in context C. When compiling separately, the effects of an **external** function are approximated conservatively by applying  $Max^r$  to the effect-annotated declared type of the function; by definition  $Max^r(\sigma)$  is  $\sigma_1 \xrightarrow[MaxEff(r)]{r} Max^r(\sigma_2)$  when  $\sigma = \sigma_1 \xrightarrow[\mu]{r} \sigma_2$ , and  $\sigma$  otherwise. This allows the view of external functions as effect-polymorphic without restricting their actual implementations.

# 6 Related Work and Conclusions

The work presented in this paper is mainly inspired by recent research on effect inference [6, 10, 11, 23, 24], efficient implementation of first-class continuations [2, 8, 22, 1], monads and modular interpreters [30, 12, 29, 13], typed intermediate languages [7, 26, 20, 17, 16, 3], and foreign function call interface [9, 18]. In the following, we briefly explain the relationship of these work with our resource-based approach.

- Effect systems. The idea of using effect-based type systems to support language interoperation was first proposed by Gifford and Lucassen [5, 6]. Along this direction, many researchers have worked on various kinds of effect systems and effect inference algorithms [10, 11, 23, 24, 28]. The main novelty of our effect system is that we imposed a "resource-based" upper-bound to the effect variables. Effect variables in all previous effect systems are always *universally* quantified without any upper bounds, so they can be instantiated into any effect expressions. Our system limits the quantification over a finite set of resources—this allows us to take advantage of the effect-resource relationship to support advanced compilation strategies.
- Efficient call/cc. Many people have worked on designing various strategies to support efficient implementation of first-class continuations [2, 8, 22, 1]. To support a reasonably efficient call/cc, compilers today mostly use "stack chunks" (a linked list of smaller stacks) [2, 8] or they simply heap allocate all activation records [22]. Both of these representations are incompatible with those used by traditional languages such as C and C++ where activation records are allocated on a sequential stack. First-class continuations thus always impose restrictions and interoperability challenges to the underlying compiler. In fact, many existing compilers choose not to support call/cc, simply because call/cc is not compatible with standard C

calling conventions. The techniques presented in this paper provide opportunities to support both efficient call/cc and interoperability with code that use sequential stacks.

- Threads. Implementing threads does not necessarily require first-class continuations but only an equivalent of one-shot continuations [1]. A finer distinction between these classes of continuations is useful, however the issues of incorporating linearity in the type system to ensure safety in the presence of one-shot continuations are beyond the scope of this paper.
- Monads and modular interpreters. The idea of using resources and effects to characterize the run-time configuration of a function is inspired by recent work on monad-based interactions and modular interpreters [30, 12, 29, 13]. Unlike in the monadic approach, our system provides a way of switching the runtime context "horizontally" from one to another via the use" (e) construct.
- Typed intermediate languages. Typed intermediate languages have received much attention lately, especially in the HOT (i.e., higher-order and typed) language community. However, recent work [7, 14, 21, 17, 3, 16, 15] has mostly focused on the theoretical foundations and general language design issues. The type system in this paper focused on the problem of compiling multiple source languages into a common typed intermediate format. We plan to incorporate the resource and effect annotations into our FLINT intermediate language [21].
- Foreign function call interface. The interoperability problem addressed in this paper has much in common with frameworks for multi-lingual programming, such as ILU, CORBA [27], and Microsoft's COM [19]. It also relates to the foreign function call interfaces in most existing compilers [9, 18]. Although these work do address many of the low-level problems, such as converting data representations between languages or passing information to remote processes, their implementations do not provide any safety guarantees (or if they do, they would require external programs run in a separate address space). The work presented in this paper focuses on interfacing programs running in the single address space with much higher performance requirements. We emphasize building a safe, efficient, and robust interface across multiple HOT languages.

We believe what we have presented in this paper is a good first-step towards a fully formal investigation on the topic of safe fine-grain language interoperations. We have concentrated on the issues of first-class continuations in this paper, but the framework presented here should also apply to handle other language features such as states, exceptions, and non-termination. The effect system described in this paper is also very general and useful for static program analysis: because it supports effect polymorphism, effect information is accurately propagated through high-order functions. This is clearly much more informative than the single one-bit (or N-bit) information seen in the simple monad-based calculus [16, 25].

There are many hard problems that must be solved in order to support a safe and fine-grained interoperation between ML and safe-C, for instance, the interactions between garbage collection and explicit memory allocation, between type-safe and unsafe language features etc. We plan to pursue these problems in the future.

#### Acknowledgment

We are grateful to the anonymous referees for their valuable comments.

#### References

- [1] C. Bruggeman, O. Waddell, and K. Dybvig. Representing control in the presence of one-shot continuations. In *Proc. ACM SIGPLAN '96 Conf. on Prog. Lang. Design and Implementation*, pages 99–107, New York, June 1996. ACM Press.
- [2] W. D. Clinger, A. H. Hartheimer, and E. M. Ost. Implementation strategies for continuations. In 1988 ACM Conference on Lisp and Functional Programming, pages 124–131, New York, June 1988. ACM Press.
- [3] A. Dimock, R. Muller, F. Turbak, and J. B. Wells. Strongly typed flow-directed representation transformations. In *Proc. 1997 ACM SIGPLAN International Conference on Functional Programming (ICFP'97)*, pages 11–24. ACM Press, June 1997.
- [4] C. Flanagan, A. Sabry, B. F. Duba, and M. Felleisen. The essence of compiling with continuations. In *Proc. ACM SIGPLAN '93 Conf. on Prog. Lang. Design and Implementation*, pages 237–247, New York, June 1993. ACM Press.
- [5] D. Gifford and J. Lucassen. Integrating functional and imperative programming. In 1986 ACM Conference on Lisp and Functional Programming, New York, August 1986. ACM Press.
- [6] D. K. Gifford et al. FX-87 reference manual. Technical Report MIT/LCS/TR-407, M.I.T. Laboratory for Computer Science, September 1987.
- [7] R. Harper and G. Morrisett. Compiling polymorphism using intensional type analysis. In Twenty-second Annual ACM Symp. on Principles of Prog. Languages, pages 130–141, New York, Jan 1995. ACM Press.
- [8] R. Hieb, R. K. Dybvig, and C. Bruggeman. Representing control in the presence of first-class continuations. In *Proc. ACM SIGPLAN '90 Conf. on Prog. Lang. Design and Implementation*, pages 66–77, New York, 1990. ACM Press.
- [9] L. Huelsbergen. A portable C interface for Standard ML of New Jersey. Technical memorandum, AT&T Bell Laboratories, Murray Hill, NJ, January 1996.
- [10] P. Jouvelot and D. K. Gifford. Reasoning about continuations with control effects. In *Proc. ACM SIGPLAN '89 Conf. on Prog. Lang. Design and Implementation*, pages 218–226. ACM Press, 1989.
- [11] P. Jouvelot and D. K. Gifford. Algebraic reconstruction of types and effects. In *Eighteenth Annual ACM Symp. on Principles of Prog. Languages*, pages 303–310, New York, Jan 1991. ACM Press.
- [12] J. Launchbury and S. Peyton Jones. Lazy functional state threads. In *Proc. ACM SIGPLAN* '94 Conf. on Prog. Lang. Design and Implementation, pages 24–35, New York, June 1994. ACM Press.
- [13] S. Liang, P. Hudak, and M. Jones. Monad transformers and modular interpreters. In *Proc.* 22rd Annual ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages, pages 333–343. ACM Press, 1995.
- [14] G. Morrisett. Compiling with Types. PhD thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, December 1995. Tech Report CMU-CS-95-226.
- [15] G. Morrisett, D. Walker, K. Crary, and N. Glew. From system F to typed assembly language. In Proc. 25rd Annual ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages, page (to appear). ACM Press, 1998.

- [16] S. Peyton Jones, J. Launchbury, M. Shields, and A. Tolmach. Bridging the gulf: a common intermediate language for ML and Haskell. In *Proc. 25rd Annual ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages*, page (to appear). ACM Press, 1998.
- [17] S. Peyton Jones and E. Meijer. Henk: a typed intermediate language. In Proc. 1997 ACM SIGPLAN Workshop on Types in Compilation, June 1997.
- [18] S. Peyton Jones, T. Nordin, and A. Reid. Green card: a foreign-language interface for Haskell. Available at http://www.dcs.gla.ac.uk:80/simonpj/green-card.ps.gz, 1997.
- [19] D. Rogerson. Inside COM: Microsoft's Component Object Model. Microsoft Press, 1997.
- [20] Z. Shao. An overview of the FLINT/ML compiler. In Proc. 1997 ACM SIGPLAN Workshop on Types in Compilation, June 1997.
- [21] Z. Shao. Typed common intermediate format. In *Proc. 1997 USENIX Conference on Domain Specific Languages*, pages 89–102, October 1997.
- [22] Z. Shao and A. W. Appel. Space-efficient closure representations. In 1994 ACM Conference on Lisp and Functional Programming, pages 150–161, New York, June 1994. ACM Press.
- [23] J.-P. Talpin and P. Jouvelot. Polymorphic type, region, and effect inference. *Journal of Functional Programming*, 2(3), 1992.
- [24] J.-P. Talpin and P. Jouvelot. The type and effect discipline. *Information and Computation*, 111(2):245–296, June 1994.
- [25] D. Tarditi. Design and Implementation of Code Optimizations for a Type-Directed Compiler for Standard ML. PhD thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, December 1996. Tech Report CMU-CS-97-108.
- [26] D. Tarditi, G. Morrisett, P. Cheng, C. Stone, R. Harper, and P. Lee. TIL: A type-directed optimizing compiler for ML. In *Proc. ACM SIGPLAN '96 Conf. on Prog. Lang. Design and Implementation*, pages 181–192. ACM Press, 1996.
- [27] The Object Management Group. The common object request broker: Architecture and specifications (CORBA). Revision 1.2., Object Management Group (OMG), Framingham, MA, December 1993.
- [28] M. Tofte and J.-P. Talpin. Implementation of the typed call-by-value λ-calculus using a stack of regions. In Proc. 21st Annual ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages, pages 188–201. ACM Press, 1994.
- [29] P. Wadler. The essence of functional programming (invited talk). In *Nineteenth Annual ACM Symp. on Principles of Prog. Languages*, New York, Jan 1992. ACM Press.
- [30] P. Wadler. How to declare an imperative (invited talk). In *International Logic Programming Symposium*, Portland, Oregon, December 1995. MIT Press.